

Enhance Reusability and Reproducibility using NCI's Provenance Capturing System

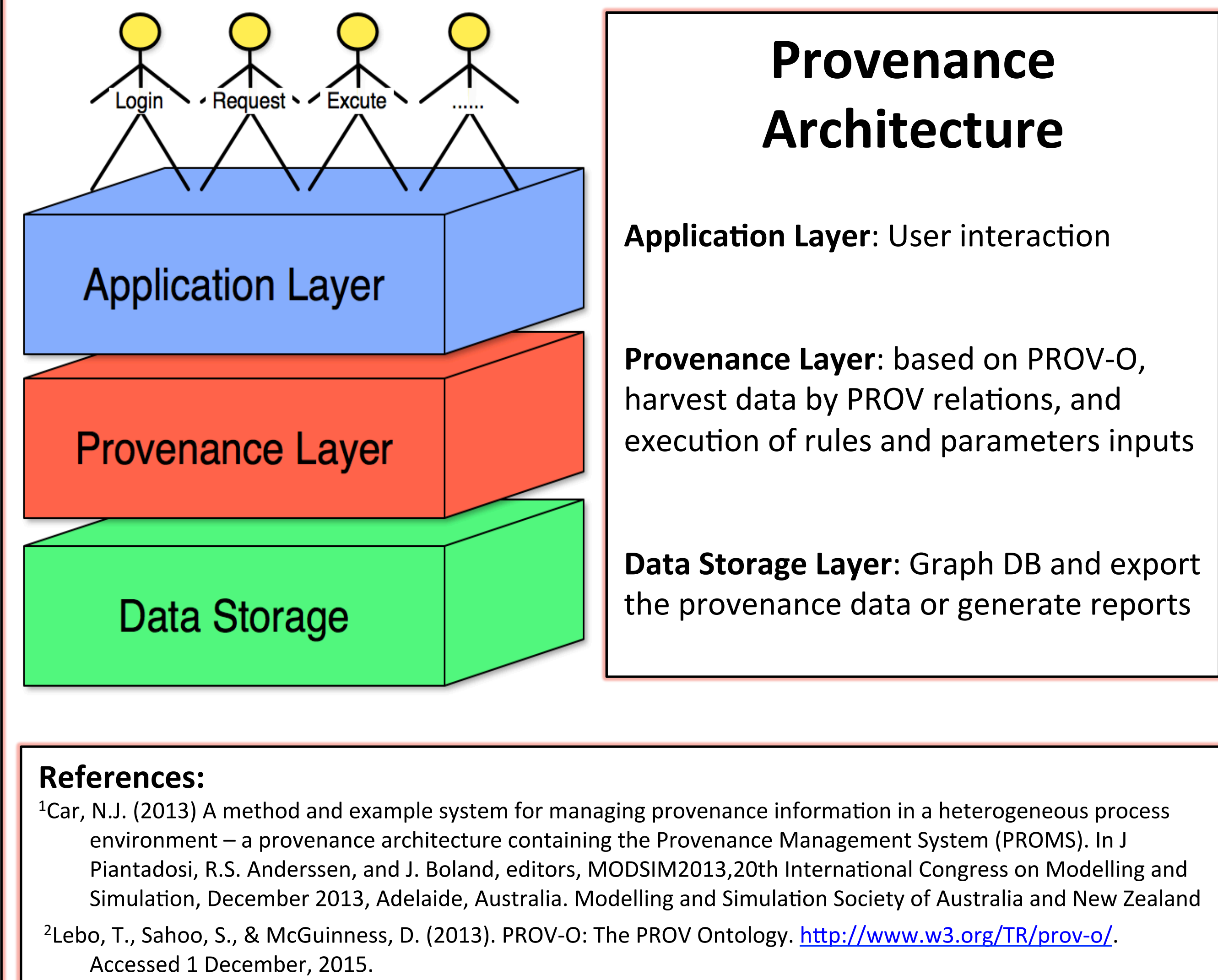
Ben Evans¹, Nick Car², Lesley Wyborn¹, Jingbo Wang¹, Wei Si¹
¹ National Computational Infrastructure, Acton, ACT; ² Geoscience Australia, Canberra, ACT

The Need for a Provenance Workflow Service

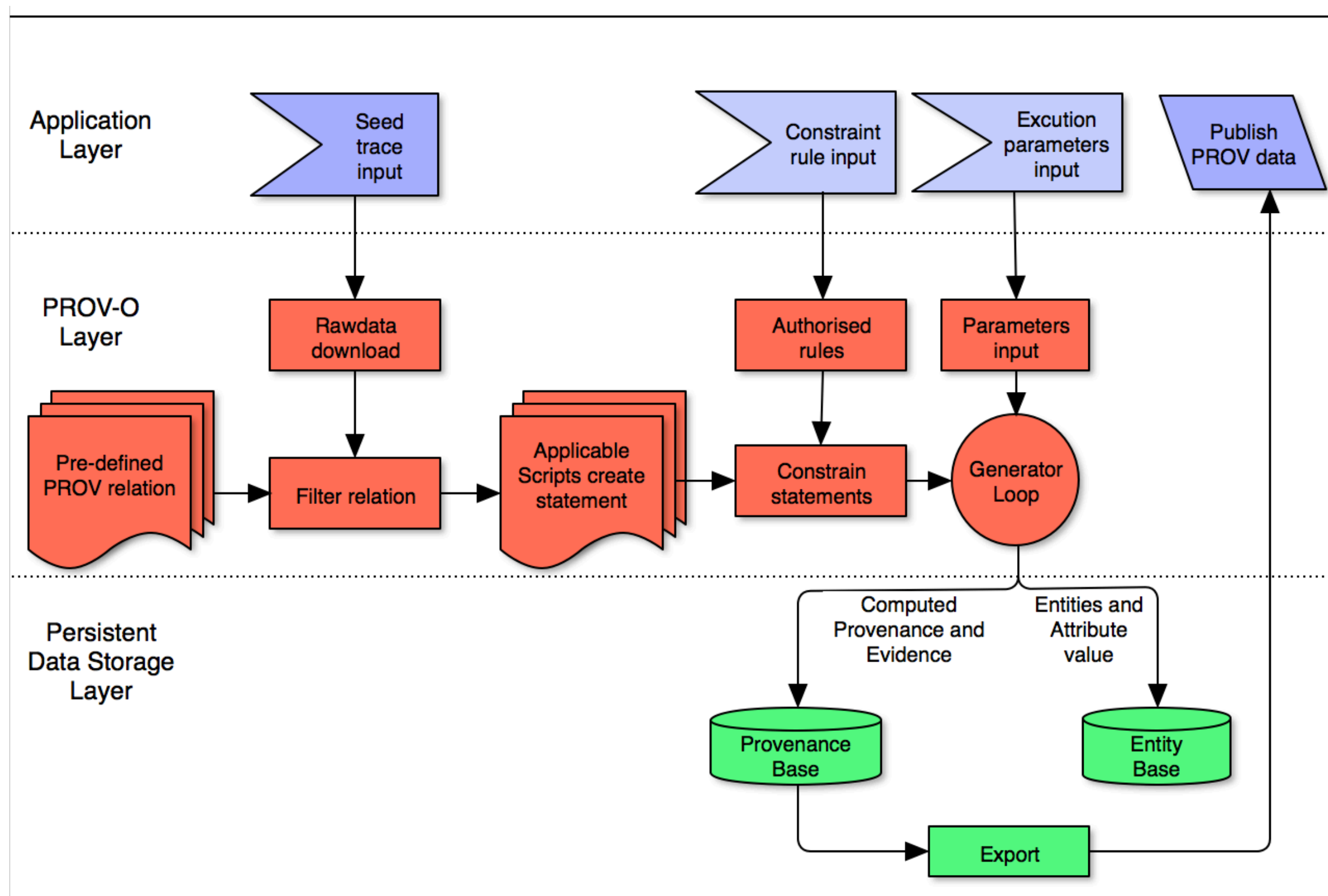
Data publication and citation have attracted considerable attention due to increasing demand to reproduce and validate scientific research outputs. Properly designed, a Provenance Workflow Service has the ability to encode transparency and accountability by providing the capability to record the key dependencies and decisions of any part of a scientific workflow, thus ensuring repeatability and trust.

NCI is using the PROvenance Management System (PROMS)¹ which provides both toolkits in a selection of programming languages for producing provenance reports compliant with PROV, the World Wide Web Consortium's provenance representation standard², and a storage system. The PROMS storage system can be queried for individual reports and elements within those reports which then link back to information in the reporting system and the data storage mechanisms it uses.

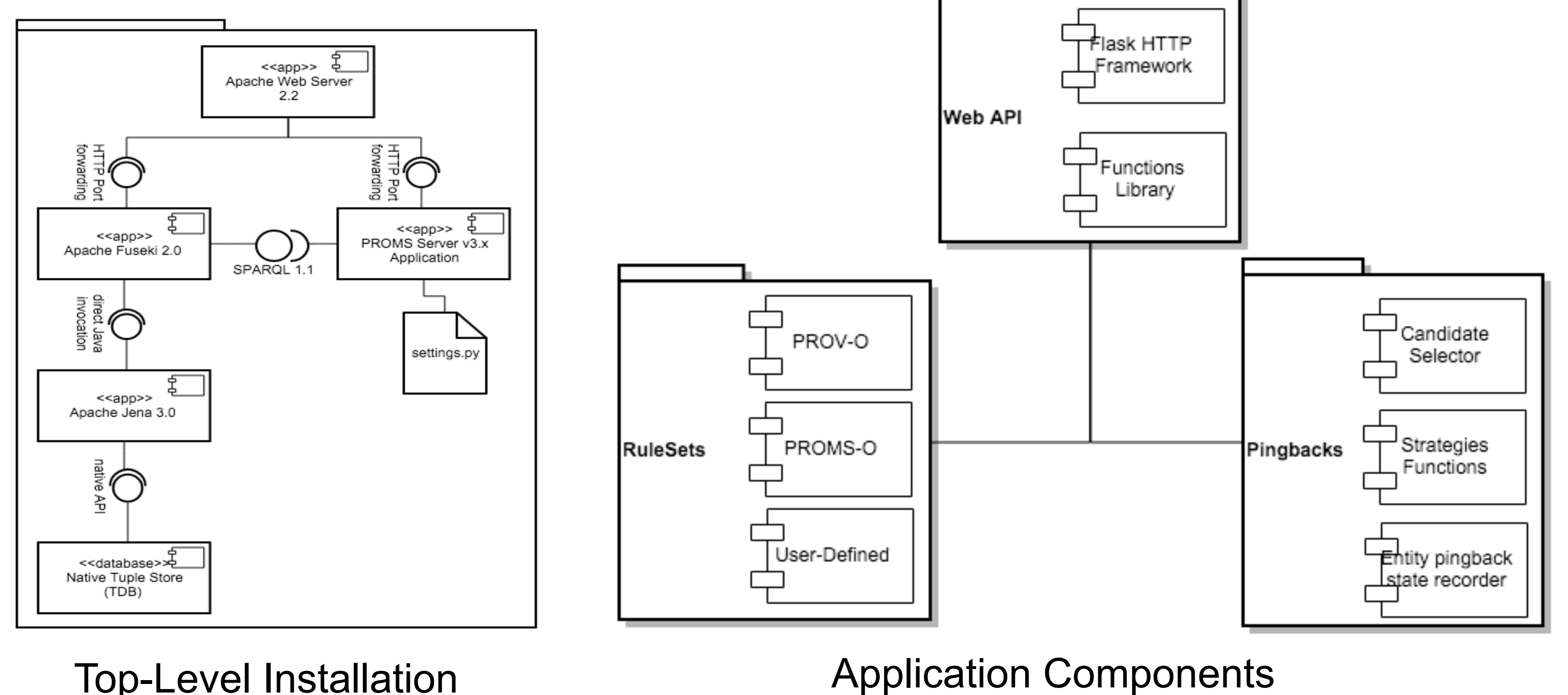
We see a well-designed comprehensive data management portal together with a provenance workflow service as an integral way of capturing the evolution of the data throughout the full data life cycle, including phases of data downloading, pre-processing and processing, re-processing and ensemble analysis. The NCI PROV server is currently in development and is available from <http://proms-dev.nci.org.au/>.



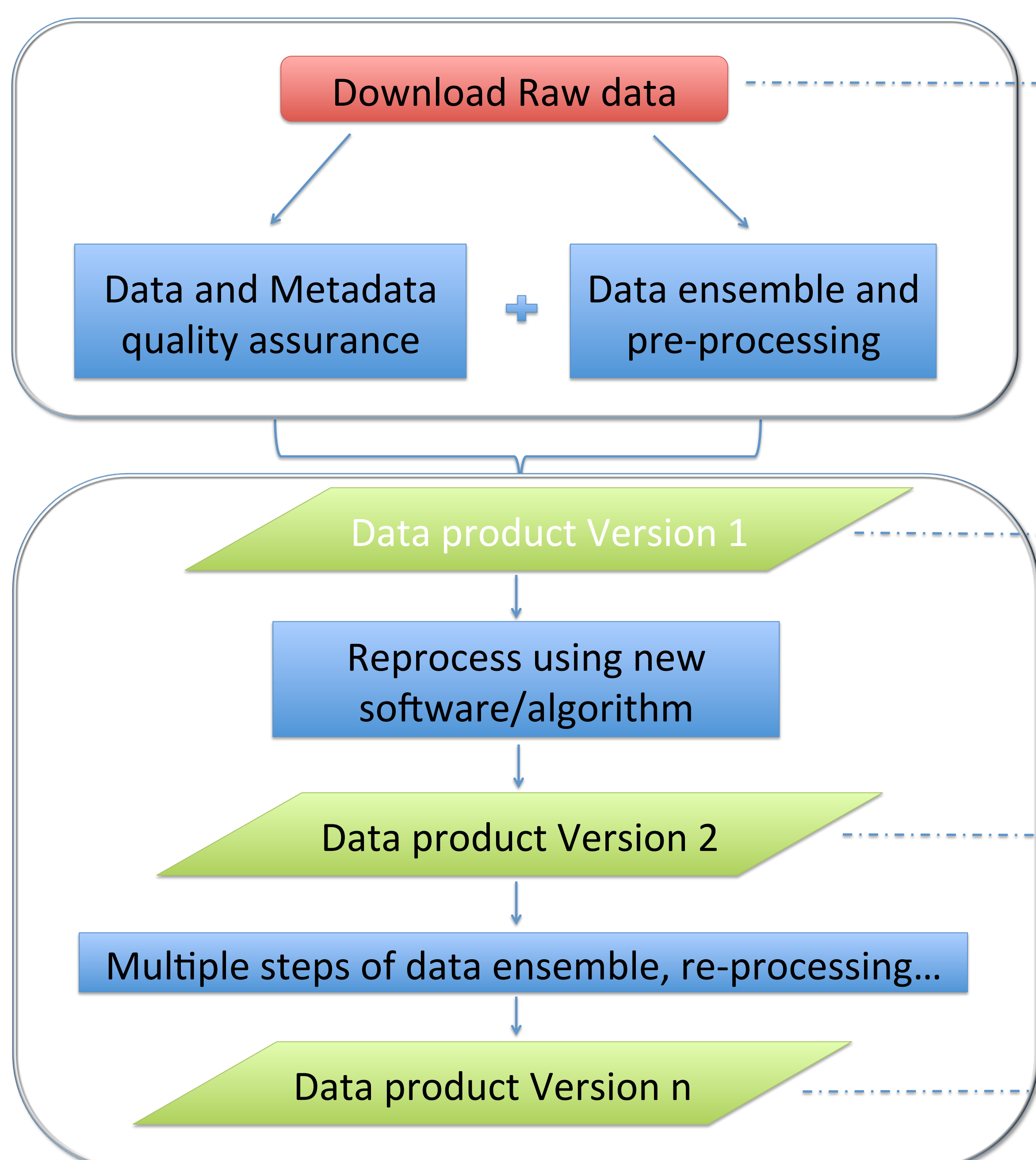
NCI's Provenance Architecture



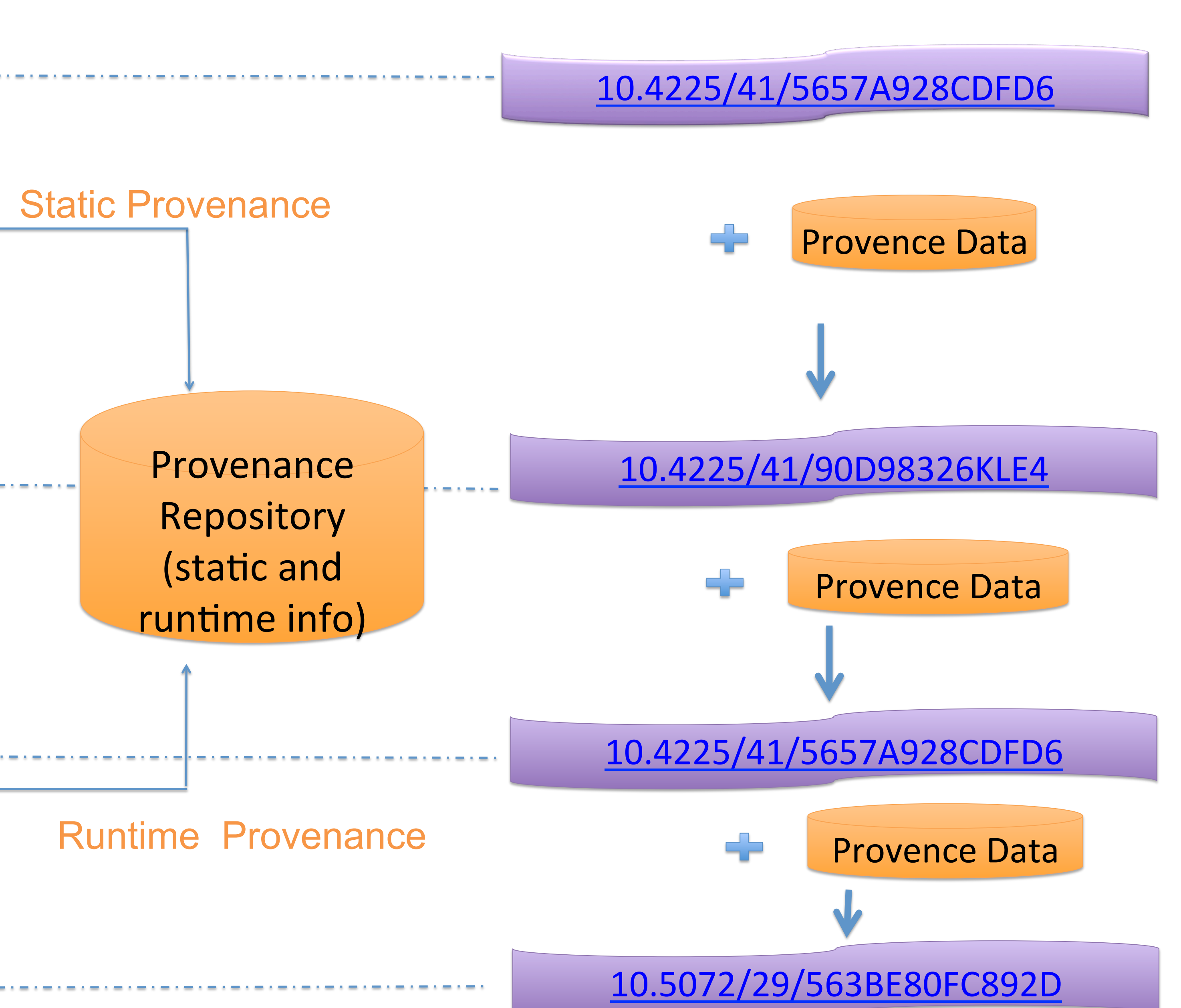
PROMS Server Architecture



Conceptual Climate Modeling Scientific Workflow



Large Scale Dynamic Citation Management



The Provenance Service captures information at each step within the end-to-end workflow, and stores it within the Provenance Repository